# Self-Supervised Representation Learning for Automatic Speech Recognition
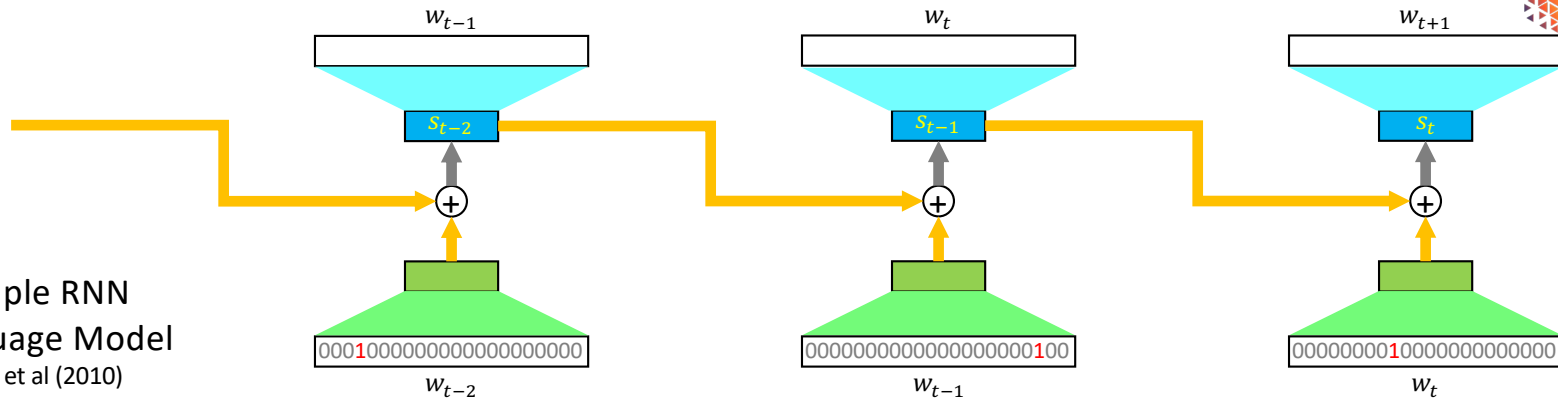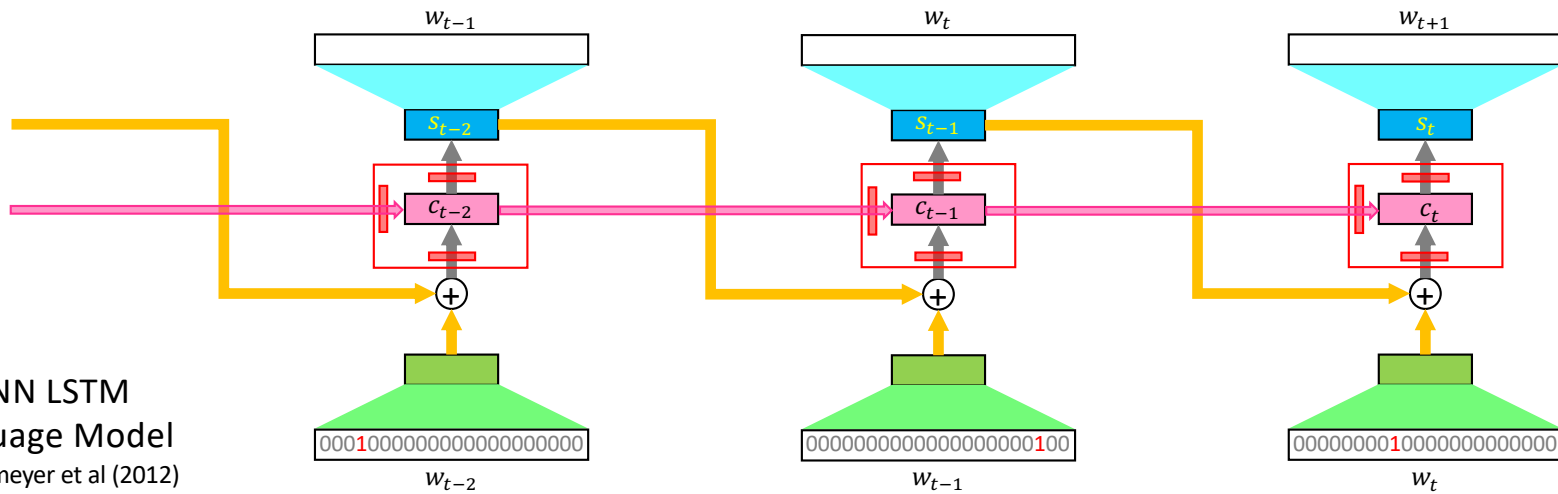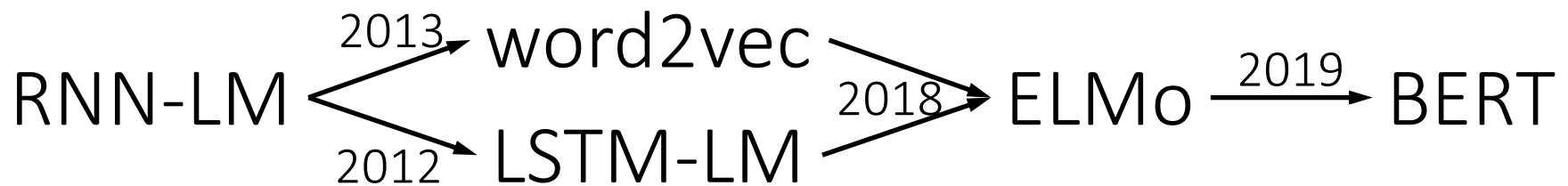
**SSL for Speech using Masked Language Model Objective: Hsu et al (2021)**
**SSL for Speech using Noise Contrastive Objective: Schneider et al (2019), Baevski et al (2020)**
**Interpreting SSL as Maximum Mutual Information Estimation**

SDAIA
الهيئة السعودية للبيانات
والذكاء الاصطناعي
Saudi Data & AI Authority

A Simple RNN Language Model
Mikolov et al (2010)

An RNN LSTM Language Model
Sundermeyer et al (2012)

RNN-LM $\xrightarrow{2013}$ word2vec $\longrightarrow$ ELMo $\xrightarrow{2019}$ BERT

$\xrightarrow{2012}$ LSTM-LM $\xrightarrow{2018}$

**Distributed Representations of Words and Phrases and their Compositionality**

**Deep contextualized word representations**

**Matthew E. Peters**[†], **Mark Neumann**[†], **Mohit Iyyer**[†], **Matt Gardner**[†],
{matthewp,markn,mohiti,mattg}@allenai.org

**Christopher Clark**[*], **Kenton Lee**[*], **Luke Zettlemoyer**[†*]
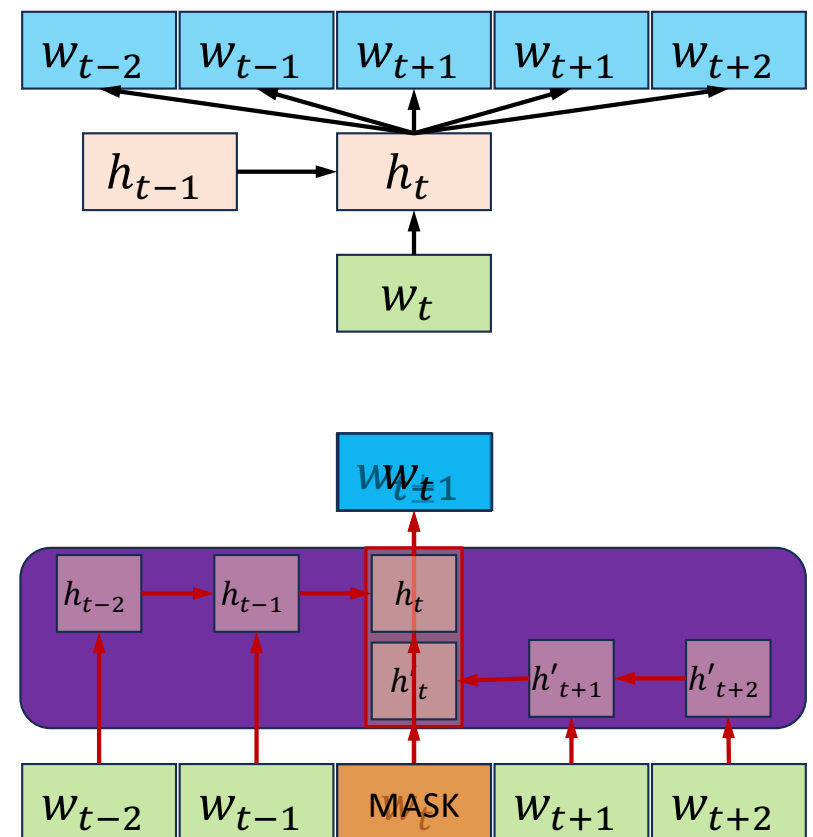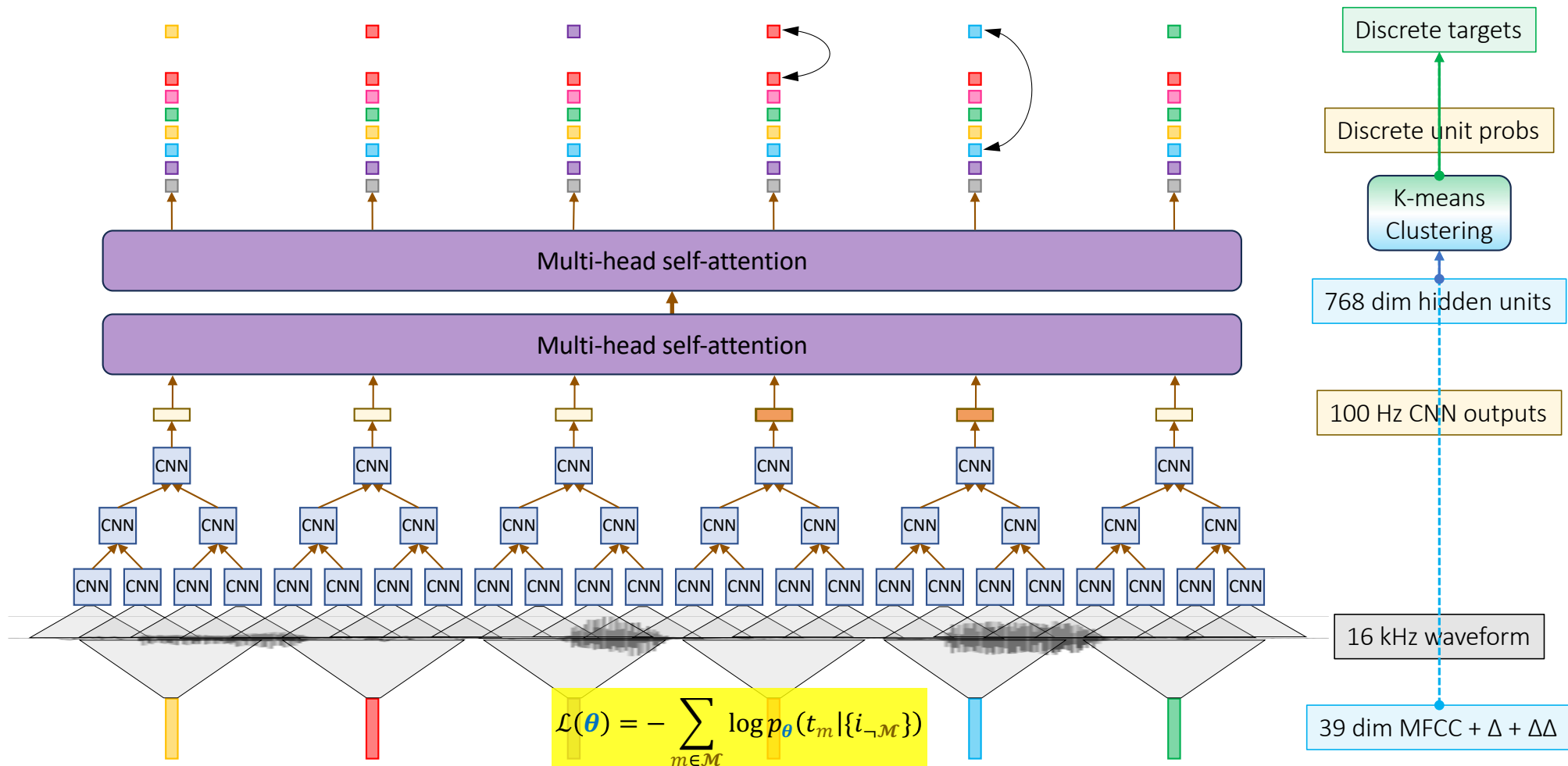{csquared,kentonl,lsz}@cs.washington.edu

**BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding**

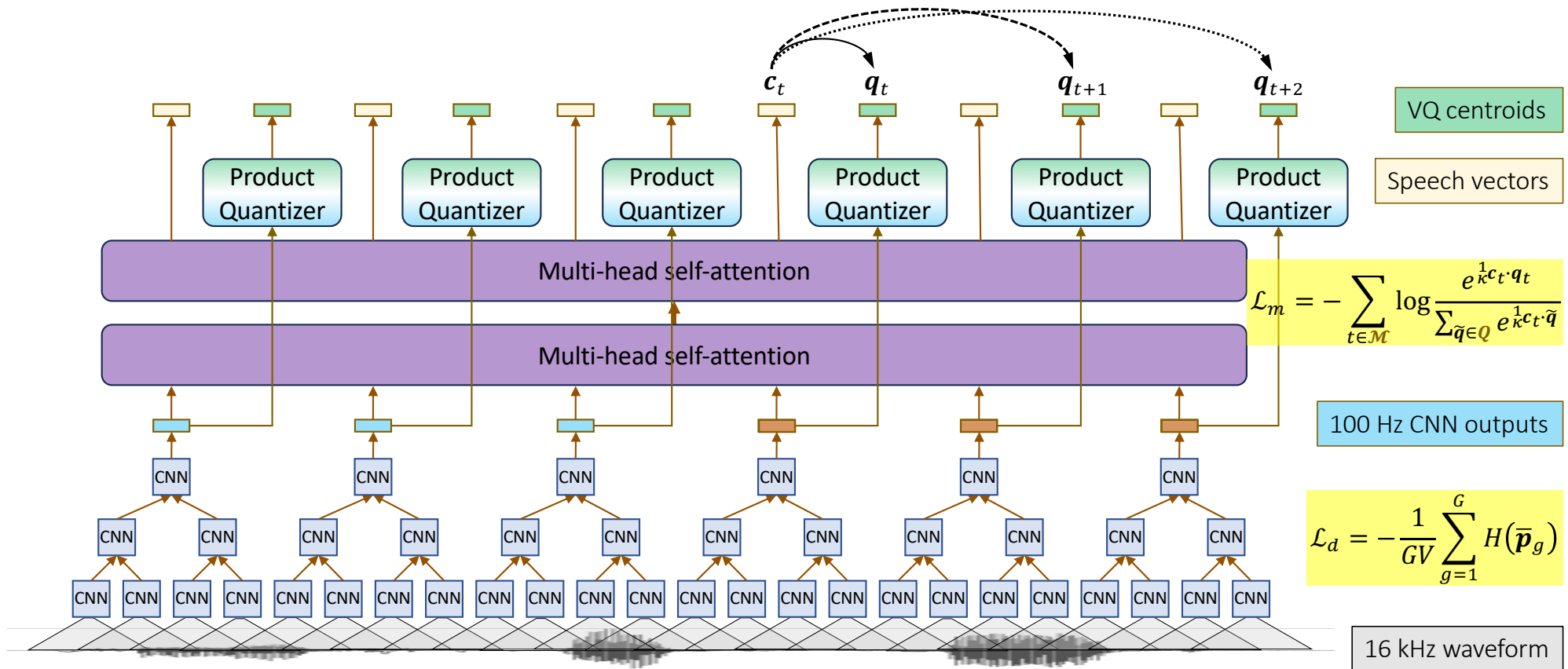**Jacob Devlin    Ming-Wei Chang    Kenton Lee    Kristina Toutanova**
Google AI Language
{jacobdevlin,mingweichang,kentonl,kristout}@google.com

| $w_{t-2}$ | $w_{t-1}$ | $w_{t+1}$ | $w_{t+1}$ | $w_{t+2}$ |

$h_{t-1} \rightarrow h_t$

$w_t$

$w_{t+1}$

$h_{t-2} \rightarrow h_{t-1} \rightarrow h_t$

$h'_t \rightarrow h'_{t+1} \leftarrow h'_{t+2}$

| $w_{t-2}$ | $w_{t-1}$ | MASK | $w_{t+1}$ | $w_{t+2}$ |

# HuBERT – Quantized speech "tokens" and BERT-like loss



$$\mathcal{L}(\boldsymbol{\theta}) = -\sum_{m \in \mathcal{M}} \log p_{\boldsymbol{\theta}}(t_m | \{i_{\neg \mathcal{M}}\})$$

Multi-head self-attention

Multi-head self-attention

CNN

16 kHz waveform

Discrete targets

Discrete unit probs

K-means Clustering

768 dim hidden units

100 Hz CNN outputs

39 dim MFCC + Δ + ΔΔ

# wav2vec 2.0 – Noise Contrastive Estimation and Learnt VQ

$c_t$    $q_t$      $q_{t+1}$      $q_{t+2}$

VQ centroids

Product Quantizer

Speech vectors

Multi-head self-attention

$$\mathcal{L}_m = -\sum_{t \in \mathcal{M}} \log \frac{e^{\frac{1}{\kappa} c_t \cdot q_t}}{\sum_{\tilde{q} \in Q} e^{\frac{1}{\kappa} c_t \cdot \tilde{q}}}$$

Multi-head self-attention

100 Hz CNN outputs

CNN

$$\mathcal{L}_d = -\frac{1}{GV} \sum_{g=1}^{G} H(\bar{p}_g)$$

16 kHz waveform

# Common misinterpretations of deep representations

Illustrated using a correctly written but often misunderstood paper

**LAYER-WISE ANALYSIS OF A SELF-SUPERVISED SPEECH REPRESENTATION MODEL**

*Ankita Pasad, Ju-Chieh Chou, Karen Livescu*

Toyota Technological Institute at Chicago
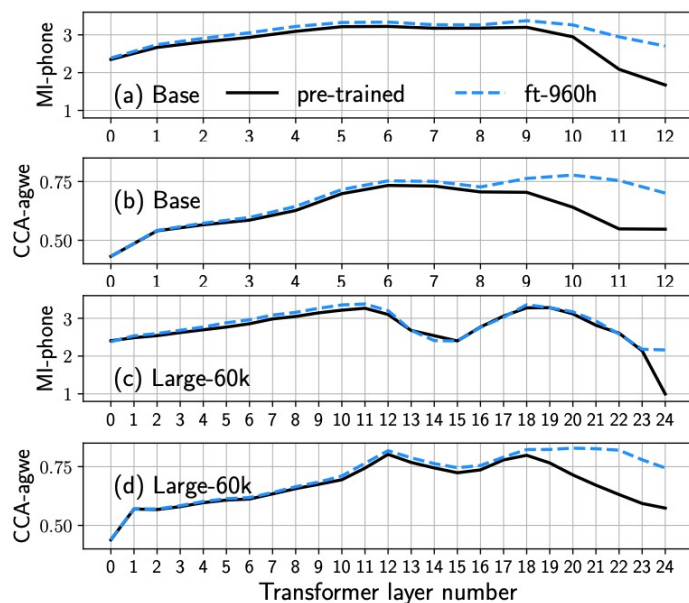{ankitap, jcchou, klivescu}@ttic.edu

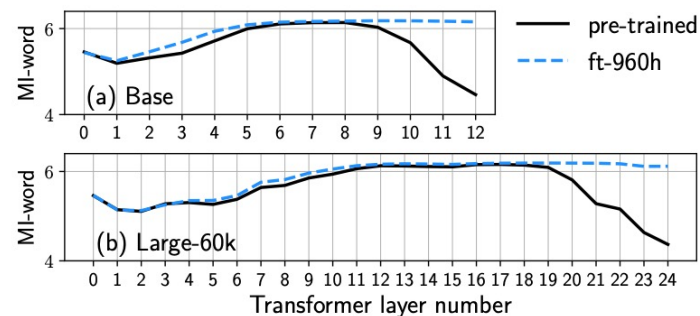Fig. 5. *MI with phone labels (max: 3.6) and CCA similarity with AGWE.*



Fig. 6. *MI with word labels (max: 6.2).*

Data processing inequality (Cover & Thomas, pp32)

$$I(A;W) \geq I(f_k(A);W) \geq I(f_l(f_k(A));W)$$

$$I(\bar{f}_{l,t_1:t_2}(A); w_i) = I(\bar{f}_{k,t_1:t_2}(A) + \bar{g}_{k:l,t_1:t_2}(A); w_i)$$
$$\gtrless I(\bar{f}_{k,t_1:t_2}(A); w_i)$$